

Hebrew Articles Category Classification

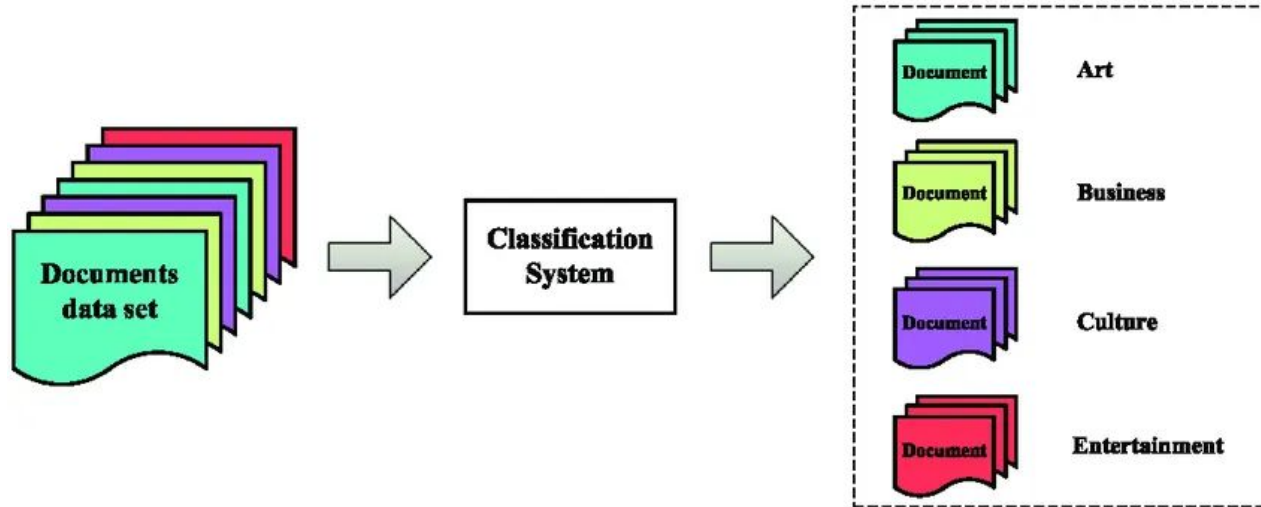
Baraa Saleh - 206466732

Marah Aboud - 208865907

A dark blue diagonal gradient bar that starts from the bottom left corner and extends towards the top right corner, covering the lower half of the slide.

Business Model

Our end goal is to create a classification model, that is able to classify articles into their respective category by leveraging data analytics and machine learning techniques

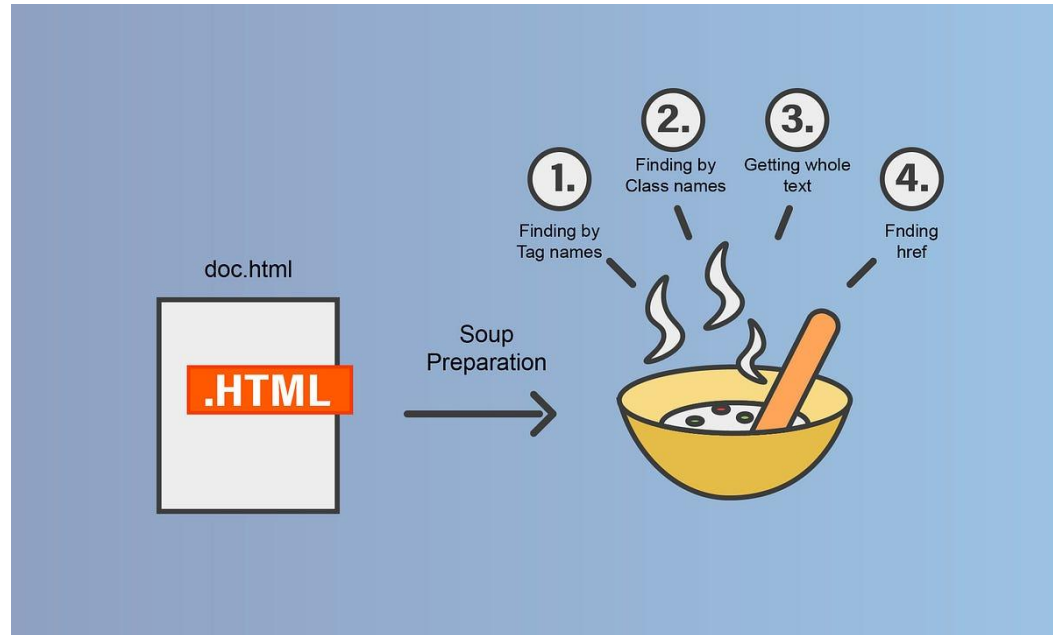


Roadmap

- Data scraping
- Preprocessing
- Classification model
- Backend: Flask
- Frontend: React

Data Scraping: Ynet + BeautifulSoup

We used Ynet which is a popular website containing articles in hebrew and parse the HTML using bs4 (BeautifulSoup) package.



Data Scraping: Headline + Description + Body

The screenshot displays the Chrome DevTools console with a JSON response from the URL `https://www.ynet.co.il/sport/article/one459w91`. The JSON structure includes fields for `headline`, `image`, `datePublished`, `author`, `publisher`, `description`, `wordCount`, `commentCount`, `genre`, `isAccessibleForFree`, `keywords`, and `articleBody`.

Key fields highlighted in red:

- `headline`: "בטיס דרמטי: 1:1 בין הפועל חדרה להפועל פ"ת"
- `description`: "לאחר 92 דקות שזוהו, טומאסביץ' כבש את הראשון עבור המארחת, אלא שאז באסי' הוכשל ברחבה, כבש מהנקודה (96) וכפה שוויון. המלאבס קרובים לירידה."
- `articleBody`: A detailed Hebrew text describing the football match between Hadera and Petach Tikva, mentioning players like Tomaszewski and Basi, and the 1-1 scoreline.

The right sidebar shows the 'Styles' panel with the 'user agent stylesheet' and 'flowplayer.css' applied to the element.

Data Scrapping: Save the data

Save the article text and its respective category into a Pandas DataFrame and then save the DataFrame into a csv.

	A	B	C	D	E	F	G	H	I	J	K
1	text	sport	entertainment	economy	health	car	food	vacation	dating	parents	environment-science
2	ההסתבכות של מנופוד הצרה סכב הפועל תא נכשל בבדיקו	1	0	0	0	0	0	0	0	0	0
3	החלוק והמסירה ער דומה מתעורר וכלי דן דוד ללדו פרננדו	1	0	0	0	0	0	0	0	0	0
4	אומצו של קין את כל הבעיות של מכבי תא אפשר היה לראו	1	0	0	0	0	0	0	0	0	0
5	שערם הפסידה 21 לפר קאסס גיא דתן הבריק דלזיה שרן	1	0	0	0	0	0	0	0	0	0
6	מתן דרזיט ובקי פודנר זג באלפות ישראל במשחה 5000	1	0	0	0	0	0	0	0	0	0
7	בידית האולסטאר באפן ברר כפי שונתב ESPN אדם סי	1	0	0	0	0	0	0	0	0	0
8	קבל אותה טרובקו נחלה ער מדלית הכפר ההיסטורית מת	1	0	0	0	0	0	0	0	0	0
9	הבתרים החדשים המאפן ברק יצחקי המנהל המקצועי אלמ	1	0	0	0	0	0	0	0	0	0
10	תעוררו ער שחר ערן זחבי סופסל ומכבי תא מצאה סקורר	1	0	0	0	0	0	0	0	0	0
11	שיר נאב אר המלחמה תשפיע ער התרבות שלו החזונית ה	0	1	0	0	0	0	0	0	0	0
12	זו לא תערכה ער השואה אלא ער חזונית הצפייה בדיכרון תש	0	1	0	0	0	0	0	0	0	0
13	מדוז נדאי לכל לחלוק נעלים בכניסה לבית או שלא רגע כל	0	1	0	0	0	0	0	0	0	0
14	מקצצים בהרש של התרבות הספרות הציבורית אינן	0	1	0	0	0	0	0	0	0	0
15	The Greater Wings של טלי ברן אלבום שמכיר בצער ארן	0	1	0	0	0	0	0	0	0	0
16	כשהטמפרטורות יורדות זו המוזיקה הישראלית שתתרו להת	0	1	0	0	0	0	0	0	0	0
17	אמני ישראל לא יכילו לעבוד בהתנדבות לצנח כולם מתקשים	0	1	0	0	0	0	0	0	0	0
18	סקוט פילצרים הולך ער זה ייתן לכל הפונה מרעשי הרקע ה	0	1	0	0	0	0	0	0	0	0
19	אוסר המתווה של השר מיני זורר לתמיכה בגופי התרבות ב	0	1	0	0	0	0	0	0	0	0
20	תמל המפנים עשנים 10 אלף הפועות צריכים עזרה כלכלית כ	0	1	0	0	0	0	0	0	0	0
21	החדש של חוליג סטונס יעקב אלף הפונה נחוצה המצלות ס	0	1	0	0	0	0	0	0	0	0
22	שנת הקאמבק של סבב נאבקו שנים כדי לא להפסיד ענשין	0	0	0	0	0	0	0	0	0	0
23	אי שורף כפר ששמיני בצד לימים קשים אבל לא חשבתי שר	0	0	1	0	0	0	0	0	0	0
24	מיליואיתמניקים שהוקרר בפעם הבאה אינשים יחשבו פעמיני	0	0	1	0	0	0	0	0	0	0
25	מחזים ארן מנכל או ממלא מקום לביטוח הלאומי ולא צפוי מ	0	0	1	0	0	0	0	0	0	0
26	טורקיה יצאה מהמשחקו כמה עולה השוללת שער בישראל כ	0	0	1	0	0	0	0	0	0	0
27	זה העסק שלי 11 שנה אבל בפעם הראשונה שנעשיתי גרפיס	0	0	1	0	0	0	0	0	0	0
28	ער רקע התקפות החותרים CMA CGM מפסיקה שוב לשוט	0	0	1	0	0	0	0	0	0	0
29	הייתי שר במסעדות מישל באה בברת לחדר אוכל ואמרה ה	0	0	1	0	0	0	0	0	0	0
30	עזיבה נוספת באוצר טון הממנה ער התקציבים חרע ער ס	0	0	1	0	0	0	0	0	0	0
31	ברושה ער שלושה ילדים ועסק חדש אני קצת פנטזיבורית וק	0	0	1	0	0	0	0	0	0	0
32	בבזבז זמן אסון הפרויקטור האזרחי סמים את תפקידו מכל	0	0	1	0	0	0	0	0	0	0
33	מחלוקת סביב מצריו המלחלה מערן בריאות המלחלה בבתי	0	0	0	1	0	0	0	0	0	0

Preprocessing: NLP

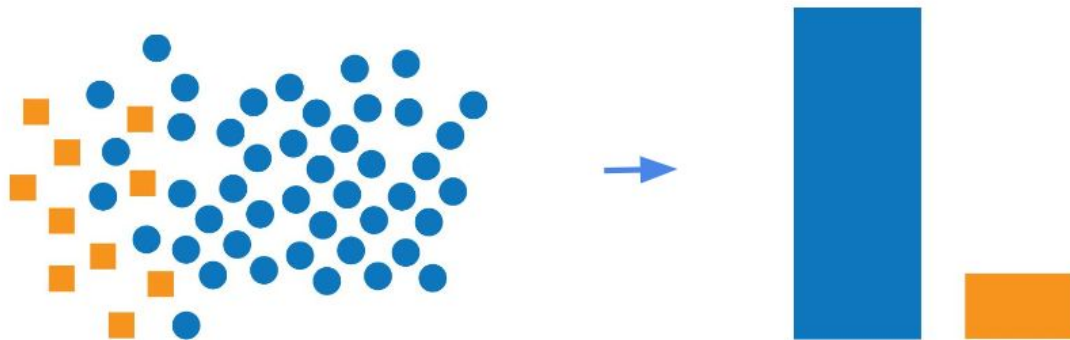
Analyze and clean the data using NLP techniques, such as removing the stop words from the texts which do not add to the understanding of the text.

For example: את, לא, של, אני, על, זה, עם, כל, הוא, אם, או, גם, ...

Preprocessing: Imbalanced Data

While checking for imbalances in the data, we noticed that 8 categories had ~900 articles and for 2 categories there were < 50 articles so we removed these 2 categories from the data due to lack of support.

Final data size ~7000 articles



Preprocessing: Encode the categories

Replace all the categories with an encoding to describe the category this text belongs to.

For example:

0 - Sport

1 - Entertainment

2 - Economy

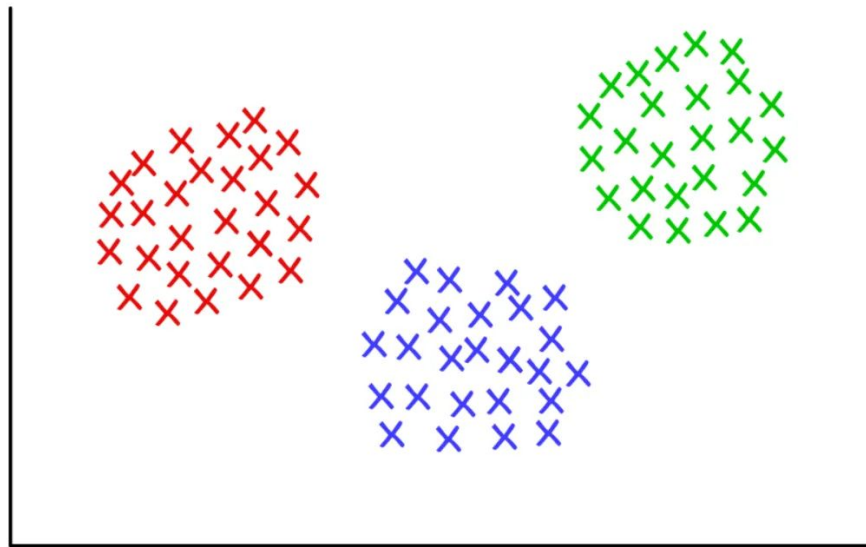
etc ...

	A	B
1	text	category_encoded
2	ההסתבכות של מנפורד הצרה כוכב הפועל תא נכשל בבריקו	0
3	החלוצ והמטרה עם רזומה מתעתע ובלו דין דוד לצדו פרנצדי	0
4	האומץ של קין את כל הבעיות של מכבי תא אפשר היה לראו	0
5	שפרעם הפסידה 21 לכפר קאסס גיא דהן הבריק דלויה שרע	0
6	מתן רוזיט ובקי פוזדנר זכו באליפות ישראל במשחה ל5000	0
7	בדיחת האולסטאר באופן ברור כפי שנכתב בESPN אדם סי	0
8	קבלו אותה גורבנקו נחתה עם מדליית הכסף ההיסטורית מק	0
9	הביתרים החדשים המאמן ברק יצחקי המנהל המקצועי אלמ	0
10	התעוררו עם שחר ערן זרבי סופסל ומכבי תא מצאה סקורר	0
11	שיר כאב איך המלחמה תשפיע על התרבות שלנו הזוועות ה	1
12	זו לא תערוכה על השואה אלא על חווית הצפייה בזיכרון הש	1
13	מדוע כדאי לכם לחלוץ נעליים בכניסה לבית או שלא רגע לפ	1
14	מקציצים בבכשת הרש של התרבות הספריות הציבוריות אינ	1
15	The Greater Wings של גולי ברן אלבום שמכיר בצער אך	1
16	כשטמפרטורות יורדות זו המוזיקה הישראלית שתרגו להת	1
17	אמני ישראל לא יוכלו לעבוד בהתנדבות לנצח כולם מתקשים	1
18	סקוט פילגרים הולך על זה ייתן לכם הפוגה מרעשי הרקע ה	1
19	אוסר המתווה של השר מיקי זוהר לתמיכה בגופי התרבות ב	1
20	חמל האמנים עשיתי 10 אלף הופעות צריכים עזרה כלכלית כ	1
21	החדש של רולינג סטונס יעניק לכם הפוגה נחוזה המלצות ס	1
22	שנת הקאמבק של טבע נאבקנו שנים כדי לא להפסיד עכשיו	2
23	אני שורף כסף ששמתי בצד לימים קשים אבל לא חשבתי שז	2
24	המיליומניקים שהופקרו בפעם הבאה אנשים יחשבו פעמיים	2
25	מהיום אין מנכל או ממלא מקום לביטוח הלאומי ולא צפוי מר	2
26	טורקיה יצאה מהמשחק כמה עולה השתלת שיער בישראל ש	2
27	זה העסק שלי 11 שנה אבל בפעם הראשונה שעשיתי גרפיט	2
28	על רקע התקופות החותמים CMA CGM מפסיקה שוב לשוט	2
29	הייתי שיי במסעדות מישלן באה גברת לחדר אוכל ואמרה ה	2
30	עזיבה נוספת באוצר סגן הממונה על התקציבים הודיע על סי	2
31	גרופה עם שלושה ילדים ועסק חדש אני קצת פנטזיזטרית וק	2
32	בבוב זמן וכסף הפרויקטור האזרחי מסיים את תפקידו מבל	2
33	מחלוקת סביב מכרז להפעלת מערך בריאות התלמיד בבתי	3

Multiclass Classification

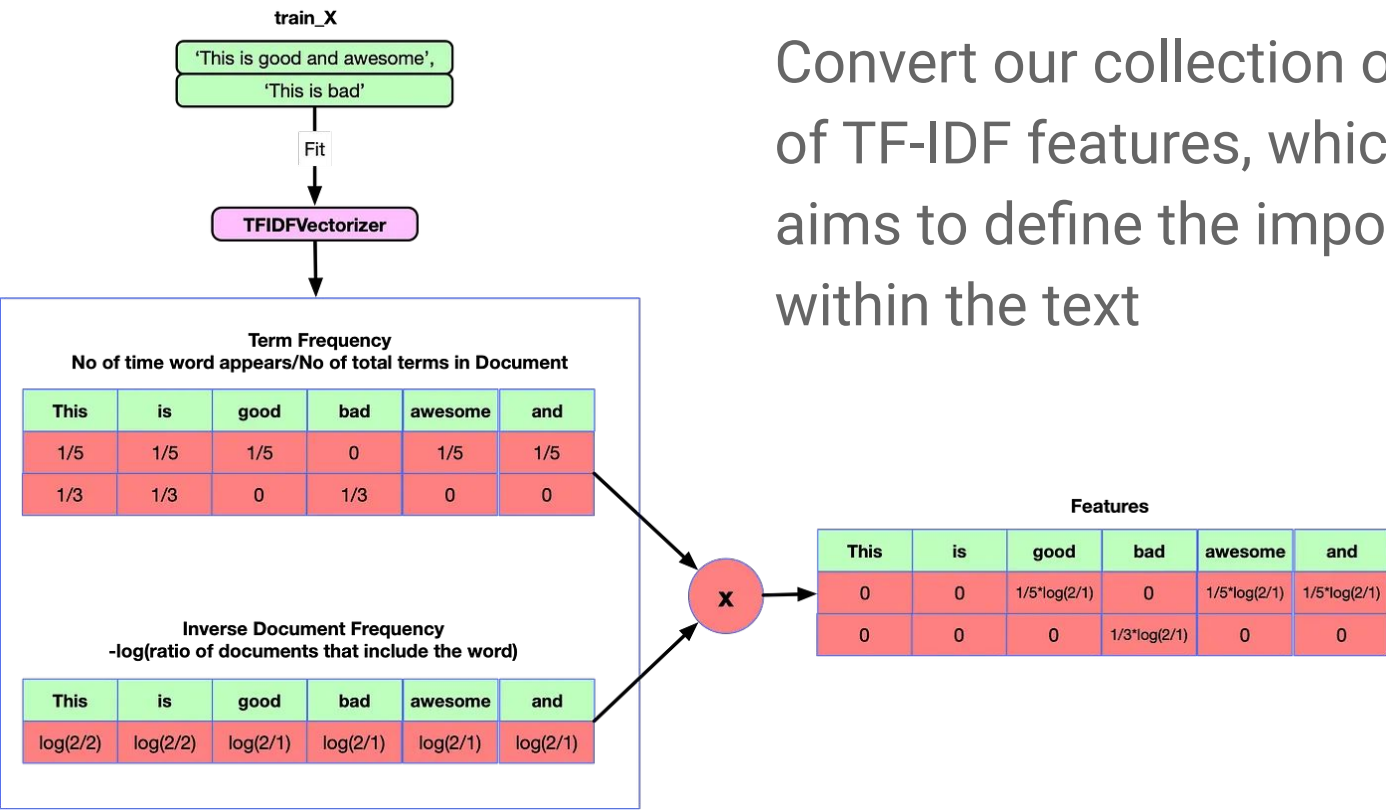
Classification is the problem of assigning observations to one or more categories.

Multiclass classification involves more than 2 categories.



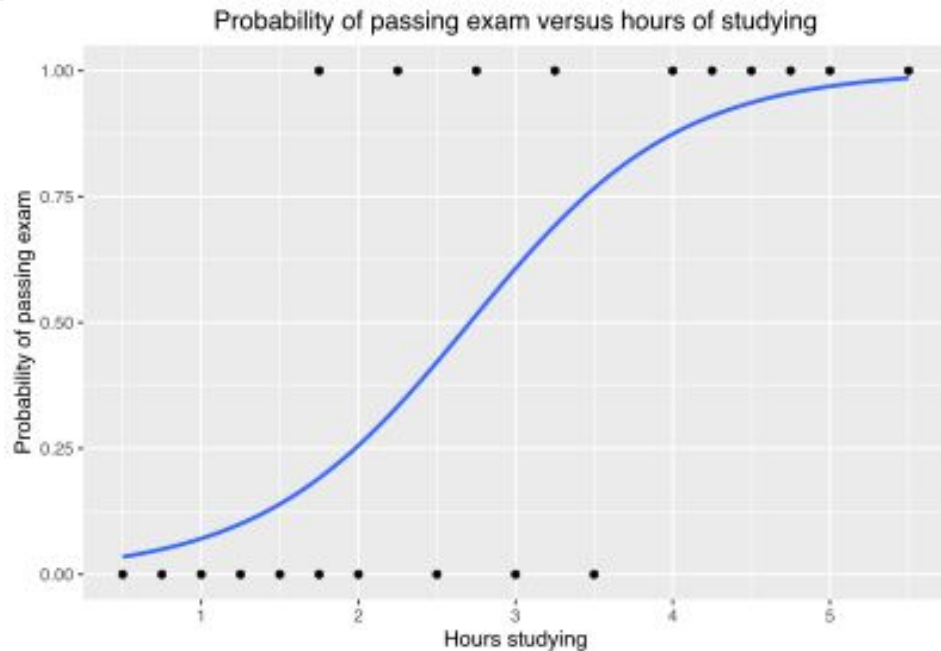
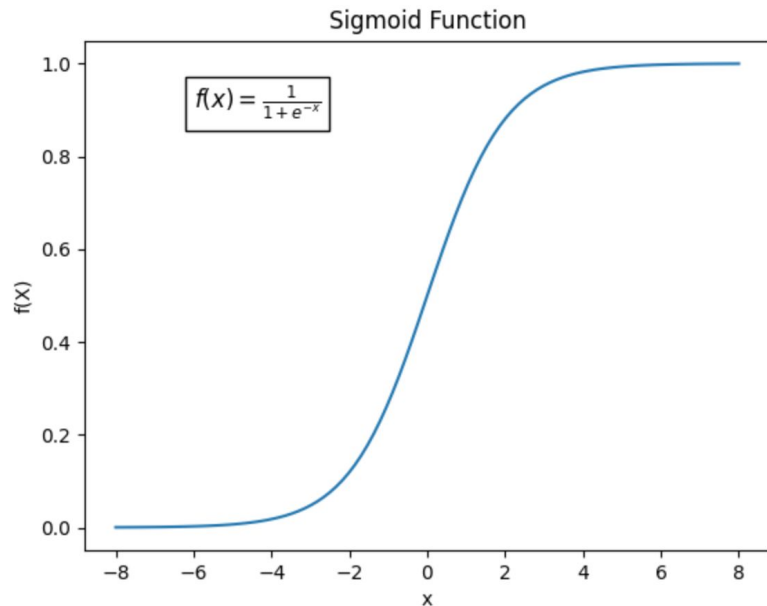
TfidfVectorizer

Convert our collection of texts to a matrix of TF-IDF features, which is a formula that aims to define the importance of a word within the text



Classification Model: LogisticRegression

LogisticRegression is a model based on sigmoid function, and here is an example of 2 classes:



Classification Model: LogisticRegression

In our case, we have multiple classes (8 to be exact) and the graph looks similar to this:

Important to note that the sum of the probabilities of a single point is 1.

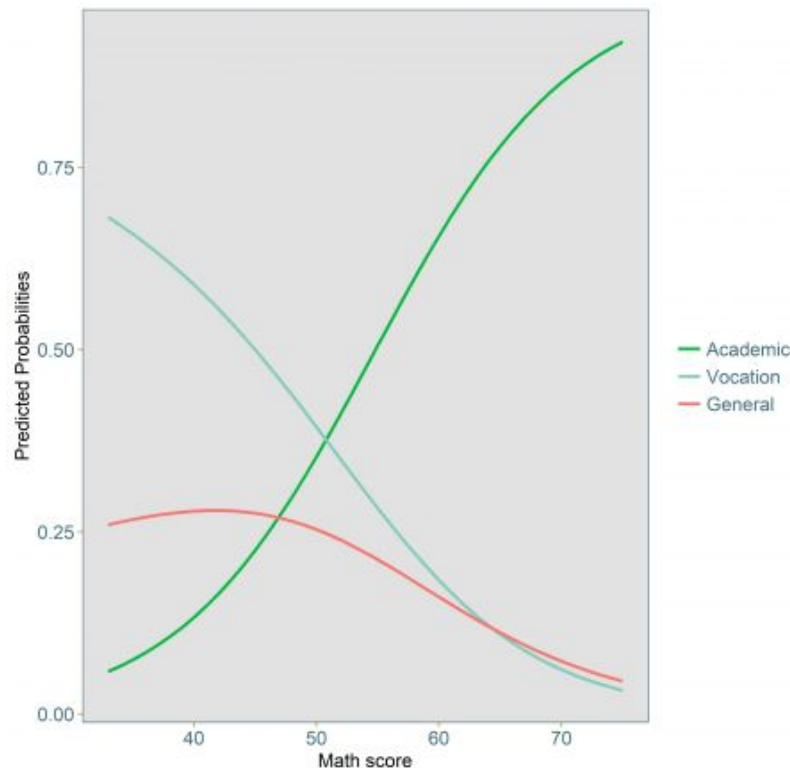
For example, if we look at grade 48, we can see that the probabilities are:

0.5 - Academic

0.25 - Vocation

0.25 - General

And their sum is 1



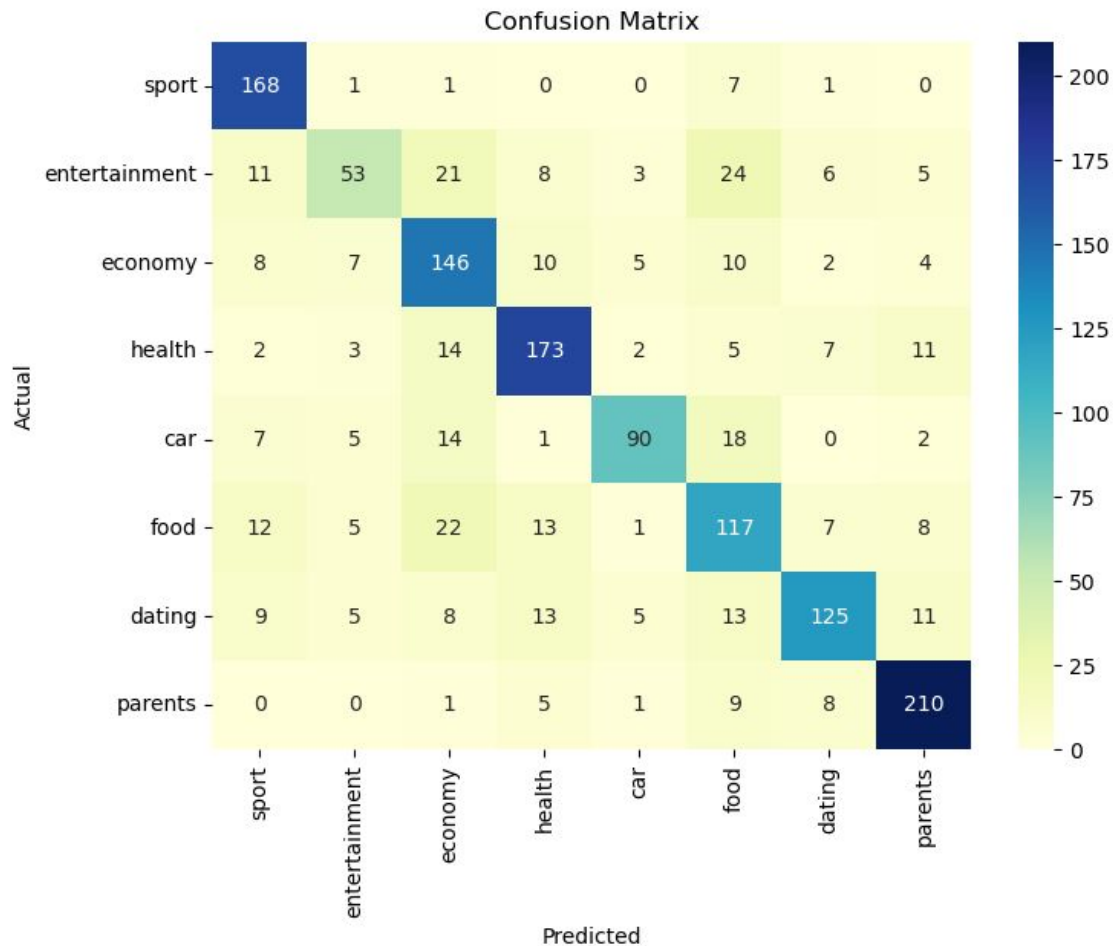
Classification Model: LogisticRegression

Our train and test process:

- Vectorize the "text" column into vectors (this is the input data X)
- Split the X and y data into X_train, X_test, y_train, y_test using train_test_split function, where the test_size is 20%
- Fit the model
- Save the model and the vectorizer into .pkl files

Confusion Matrix

Testing the model on %20 of the data and see how many times the model answered correctly



Model Performance Metrics

LogisticRegression				
	precision	recall	f1-score	support
0	0.77	0.94	0.85	178
1	0.67	0.40	0.50	131
2	0.64	0.76	0.70	192
3	0.78	0.80	0.79	217
4	0.84	0.66	0.74	137
5	0.58	0.63	0.60	185
6	0.80	0.66	0.72	189
7	0.84	0.90	0.87	234
accuracy			0.74	1463
macro avg	0.74	0.72	0.72	1463
weighted avg	0.74	0.74	0.73	1463

Model Performance Metrics

Here are a few metrics that helps check model performance.

For example, let's look at the Sports category:

$$\text{Precision} = 168 / (168 + 49) = \underline{\mathbf{0.77}}$$

$$\text{Recall} = 168 / (168 + 10) = \underline{\mathbf{0.94}}$$

$$F1 = (2 * 0.77 * 0.94) / (0.77 + 0.94) = \underline{\mathbf{0.85}}$$

$$\text{Accuracy for all} = 1082 / 1463 = \underline{\mathbf{0.74}}$$

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$F_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

Backend + Frontend

Backend: Flask server loads the model the vectorizer and listens for requests.
Frontend: React framework, we send requests to the backend and show the results in a table.

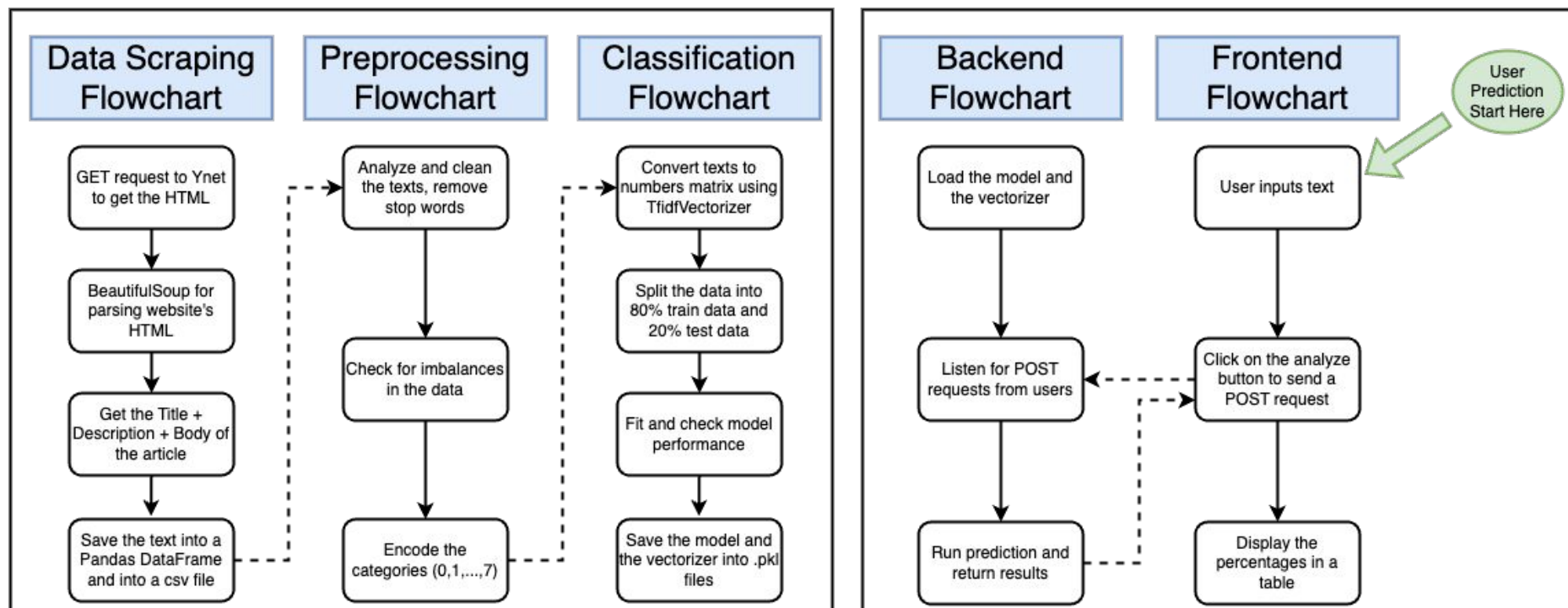
סיווג טקסטים לפי קטגוריות

X	Headers	Payload	Preview	Response	Initiator	Timing
▼	{0: {confidence_percentage: 0.8927070539192865, english: "Sport", hebrew: "ספורט"}, 1: {confidence_percentage: 0.016499586357066467, english: "Entertainment", hebrew: "תרבות"}, 2: {confidence_percentage: 0.017083240317317585, english: "Economy", hebrew: "כלכלה"}, 3: {confidence_percentage: 0.01433074011155041, english: "Health", hebrew: "בריאות"}, 4: {confidence_percentage: 0.01738512554484206, english: "Car", hebrew: "רכב"}, 5: {confidence_percentage: 0.020247947447207933, english: "Food", hebrew: "אוכל"}, 6: {confidence_percentage: 0.013065666603139063, english: "Dating", hebrew: "יחסים"}, 7: {confidence_percentage: 0.008680639699589945, english: "Parents", hebrew: "הורים"}}					

קטגוריה - Category	אחוז דואות ↓
ספורט - Sport	89.271 %
אוכל - Food	2.025 %
רכב - Car	1.739 %
כלכלה - Economy	1.708 %
תרבות - Entertainment	1.65 %
בריאות - Health	1.433 %
יחסים - Dating	1.307 %
הורים - Parents	0.868 %

למורסיה: מרקו תודורוביץ' 23 נקודות, דסטין סליבה 13 נק', דילאן אניס 10 נק', ג'ונה ראדנוב 9 נק', מוסה דיאן ורודוניוס קורוקס 2 נק' כ"א וטרי קופיין נקודה אחת. רבע ראשון - 16:17
למורסיה: חמישיית הפועל חולון: ניב משגב, דרו קרופורד, סי ג'יי האריס, קווין הרווי וג'סטין סמית'. חמישיית מורסיה: טרי קופיין, דילאן אניס, רודוניוס קורוקס, דסטין סליבה, ומרקו תודורוביץ'. סליבה וקורוקס היו אחראים לחמש הנקודות הראשונות של הספרדים, כשהרווי היה החזיר מיד עם אחת משלו. ניב משגב והרווי קלעו מחצי מרחק ושלשה של האריס העלה את החבורה של עמית שרף ליתרון ראשון בהתמודדות. היתרון החל להתנדנד מצד לצד כשג'ונה ראדוב קלע מעבר לקשת מהפינה ובצד השני טאג'יר מקול דייק פעמיים מהקו. הגארד של חולון נחסם במהלך הראשון של הרבע וילוח התוצאות הראה 16:17 למורסיה לאחר 10 דקות של כדורסל רבע שני - 31:33
למורסיה: שחר עמיר קלע שלשה במהלך הראשון של הרבע, מקול קלע גם כן בפתיחה כשבצד השני ארבע נקודות של ראדוב קבעו 21:21. האנס ואגוניו הוסיפו שלשה ושתי נקודות של הרווי מהקו העלו את הסגולים ליתרון חמש. ההייליט הגדול של הערב הגיע מניב משגב שהריס האלי הופ מחצי מגרש לג'סטין סמית' שסיים בדאנק. סיטו אולנסן עצר את הדריירה של חולון עם פסק זמן. סמית' עם סל ועבירה העלה את ההפרש לשמונה. תודורוביץ' קלע והוסיף נקודה מהקו ונקודות של אניס צימקו את ההפרש לשלוש בלבד. תודורוביץ' וסליבה השלימו ריצת 0:10 שהחזירה את היתרון לספרדים וגרמה לעמית שרף להזעיק פסק זמן. סמית' ומקול החטיאו זריקות פנויות מתחת הסל לקראת סיום הרבע ובסיומו הישראליים ירדו בפיגור 33:31 לאחר שלא קלעו נקודות בכמעט ארבע דקותיו האחרונות. רבע שלישי - 43:56 למורסיה: הלמרות נקודות של ניב משגב, פתיחת הרבע הייתה של הספרדים ושלשה של תודורוביץ' העלתה את היתרון שלהם לחמש הפרש. הסגולים ניצלו בזכות כמה הטאות קלות

Flowchart



Technologies

- Data scraping: BeautifulSoup, Pandas
- Preprocessing: NLP, Imbalanced data, Encoding
- Classification model: TfidfVectorizer, LogisticRegression
- Backend: Flask
- Frontend: React

סיווג טקסטים לפי קטגוריות

קטגוריה - Category	אחוז ודאות ↓
Sport - ספורט	0 %
Entertainment - תרבות	0 %
Economy - כלכלה	0 %
Health - בריאות	0 %
Car - רכב	0 %
Food - אוכל	0 %
Dating - יחסים	0 %
Parents - הורים	0 %

יש להזין את הטקסט...

פענח